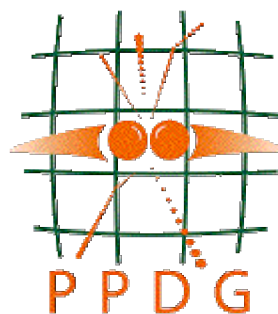


Particle Physics Data Grid Collaboratory Pilot

**PPDG Cross-Project Summary**  
**PPDG Steering Committee,**  
 May 29<sup>th</sup> . 2004



Introduction .....	1	US CMS .....	5
STAR .....	1	BaBar .....	6
JLab .....	2	Condor Group .....	7
D0 .....	3	Globus Team .....	7
LHC Experiments .....	4	Storage Resource Broker .....	8
US ATLAS .....	4	Storage Resource Manager .....	8

## Introduction

PPDG conducts an annual cross-project document or review process. For 2004, the third year of the SciDAC collaboratory pilot, we present an extended version of the Science Impact note that was presented around the time of the SciDAC PI meeting in March 2004. This document does not include work in progress, or address the needs and issues of the project.

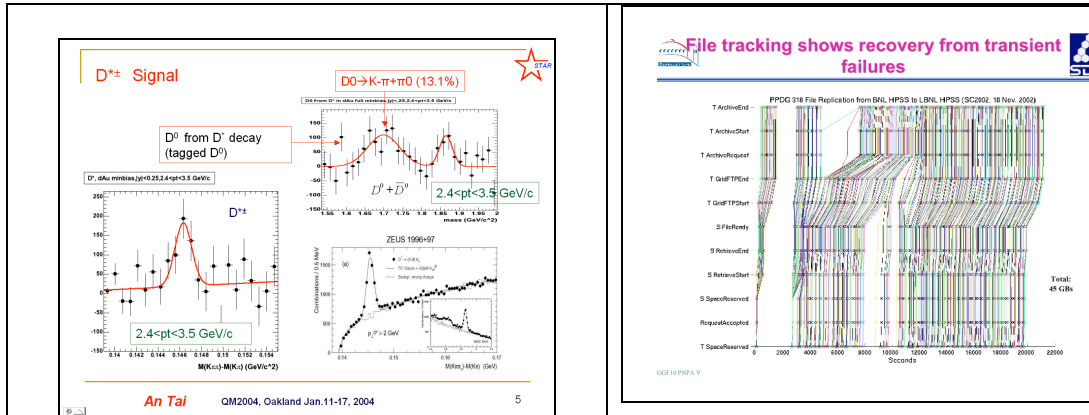
Previous deliverables in this series have been: Experiment Requirements (part of the proposal text), 2001; Sub-project reviews, 2002 ([PPDG-15](#)); Experiment Questionnaires and responses, 2003 ([PPDG-30](#))

## STAR

A wealth of physics results have emerged from the RHIC program and STAR has published 26 papers in refereed journals, 12 of which were published since the start of 2003. The first publication from the d+Au run in the spring of 2003 was in print in a record four months after the completion of data taking. This paper, along with those of the three other RHIC experiments and previous results, provided strong evidence that the density of the medium generated in central Au+Au collisions is many times that of cold nuclear matter. An example of the first direct measurement of open charm production at RHIC is shown below and provides the first salvo of a full program of measurements utilizing the large datasets recorded in the current and future runs at RHIC.

The STAR computing groups at BNL and LBNL achieve sustained and robust automated data transfers of 5 TB a week in both directions between the HPSS mass storage systems at BNL and LBNL. This allows “next day” access to fresh data for analysis and physicists using the facilities at LBNL and BNL are able to collaborate more effectively on the analyses that have led to the recent physics results. Performance is measured regularly and monitoring is used to spot and identify problems in the data transfer. The figure below illustrates the robustness of the end-to-end system even though the time trace for individual files shows delays due to various component

errors or down times.



Traditionally transferring 100's of Gigabytes of data and 1000's of files was tedious, time consuming, error prone and insecure. Overall the process would take up to 10 days for < 1 TB of data and would take a large fraction of an FTE, and still result in about 1% inconsistency in the number of files in a dataset between the two sites. Transfer of TB datasets are now accomplished using grid technologies with a single command and the average throughput is 10 times greater mostly due to improved operational efficiencies. Unsuccessful transfers due to transient problems in the end-to-end system, though rare, may be corrected easily by comparing the BNL and NERSC file catalogs to generate a secondary list of files to be transferred. A file discrepancy rate of 0.02% or less is now obtained, 50 times less than before. Reliable file transfer and catalog registration not only benefits the scientists and the experiment's physics throughput but also considerably reduces the tedious everyday task of the operations staff.

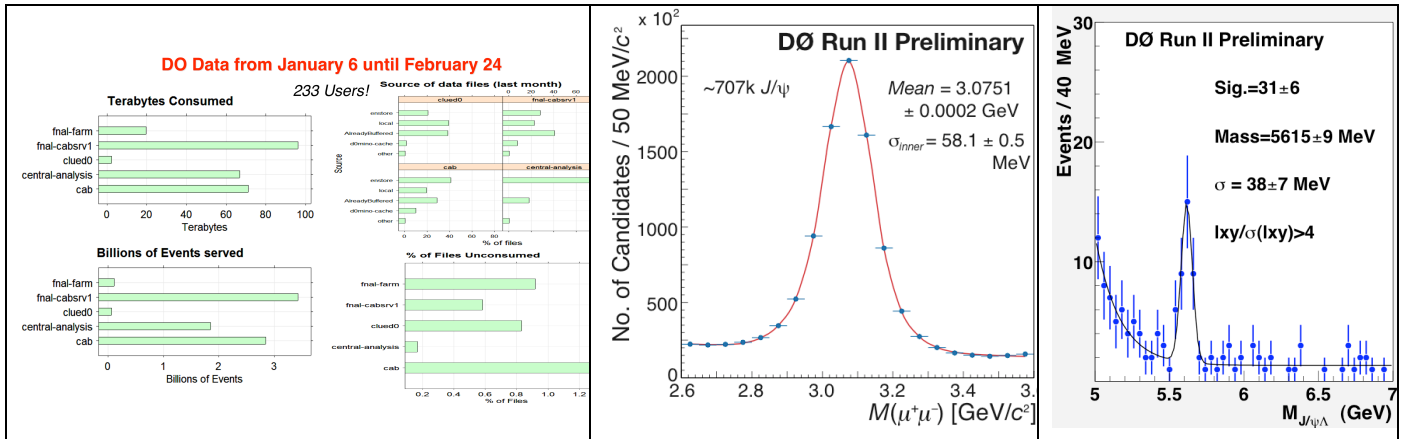
References: An Tai, "STAR Measurements of Open Charm", Quark Matter 2004, <http://qm2004.lbl.gov/>; Monitored transfers: <http://sdm.lbl.gov/datagrid/ppdg/2004/>; Storage Resource Management <http://sdm.lbl.gov/srm>.

## JLab

The experiment groups at FSU, Glasgow, and ODU have used the SRM service to transport data between their site and JLab in order to achieve more rapid completion of both simulations and data analysis. The time to simulate 30M events was reduced from 3+ months to 1 week, a 10X reduction in time, because of ability to effectively use FSU facilities and transport the data back to the JLab Mass Storage System. Based on these trials, the FSU group anticipates being able to reduce time delays for calibrations after data is acquired by CLAS once SRM V2 goes into production this year. Data transport across the Atlantic between JLab and Glasgow underwent a paradigm shift from tapes by FedEx to network data transport when the SRM service was provided at JLab. The result was a savings of both time and money. Without the SRM service, network data transports required a lot of interaction and time. The SRM automated the data staging process and data transfers.

## D0

D0 developed has been steadily using and further developing its distributed data handling system Sequential Access using Metadata (SAM) for simulation montecarlo data creation and processing before the start of data taking, and for data acquisition, processing and analysis of data acquired from the Tevatron over the last two years. During the last six months a reprocessing of the complete dataset has been done. 20% of the 500 million events of reprocessing was done off-site from Fermilab, with the input raw and output processed event data being trasmitted using GridFTP. Without the availability of this grid technology the reprocessing and presented physics results could not have been completed in time for the Spring conferences at Morion in March. Over 50 Terabytes of event data were transported with GridFTP into the central mass storage system at Fermilab for reprocessing. 10s of Terabytes have been transmitted for montecarlo and offsite analysis (there is of the order of 500 Terabytes stored in the mass storage system for D0). Using multiple streams in GridFTP has increased the throughput rate a factor of 5 over previous FTP technology.



Through PPDG effort, Grid based job scheduling, management and monitoring (JIM) has been introduced into the system. Because production processing and analysis of the D0 data is in full swing great care must be taken to preserve the robustness of the overall system and introduce these new capabilities slowly. The Grid job and information management technologies – Condor-G, Globus Gatekeeper, Dagman, MDS – have to be integrated in all details with the existing data handling and job execution environments.

In 2004, more than 8000 jobs have been run over 3 sites using the JIM based job planning and management tools based on Condor-G matchmaking (for which extensions were developed as part of PPDG) and job execution. Early use of the JIM technology have convinced the physicists that gains in effort of 1-2 FTEs over the 5 FTEs needed for the recent reprocessing are achievable in the next 12 months and that the efficiency of processing results will be of the order of a factor of 50%.

References: SAM <http://d0db.fnal.gov/sam/> ; reprocessing: <http://www-d0.fnal.gov/computing/reprocessing/> ; D0 collaboration meeting: [http://d0server1.fnal.gov/projects/meetings/collab\\_feb\\_2004/collab\\_feb\\_2004.html](http://d0server1.fnal.gov/projects/meetings/collab_feb_2004/collab_feb_2004.html)

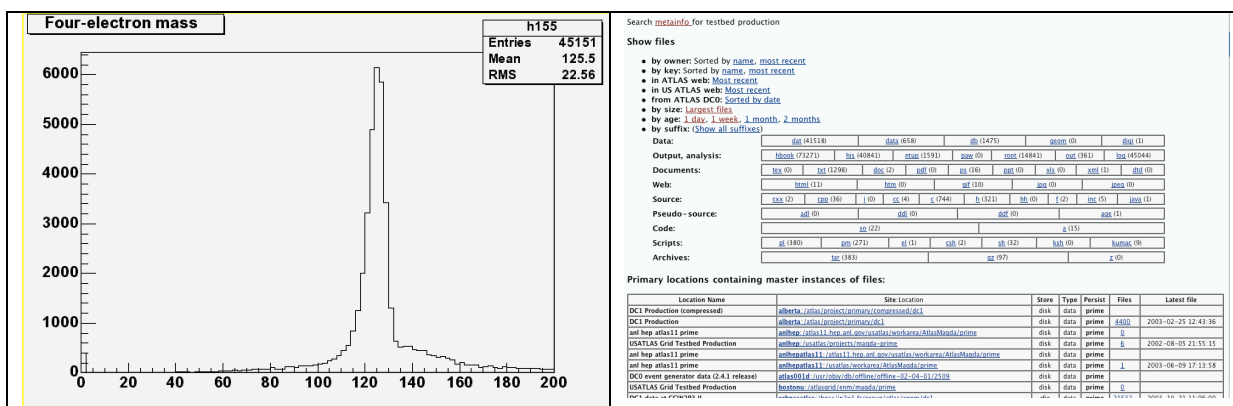
## LHC Experiments

By 2007, ATLAS and CMS must have in production, operationally supported and proven, robust and performant data analysis and processing systems to provide the necessary reconstruction and analysis of events of about 1 MB in size archived at a rate of 100HZ. Well before that time the experiments must have sufficient statistically generated and subsequently analysed hundreds of millions of simulated events to enable decisions to be taken for the physics triggers and detector commissioning, as well as development of appropriate computing and physics models. The collaborations meet these challenges through a series of annual Data Challenges of increasing scale and complexity, with specific milestones scheduled to meet the needs of the collaboration wide experiment detector and physics analysis groups. Benefits from grid based systems come through the increased number of events generated, the efficiencies in peoples effort through automation, and the trust in and ability to validate the generated and stored data.

## US ATLAS

US ATLAS DC1 metrics (in 2002-2003) were to generate and record about 2% of the event rate expected at experiment turn on. The MAGDA system developed through PPDG was used in the U.S. for data transfer and recording. Over 40 TB have been recorded in MAGDA from the distributed sites in the U.S, and an additional 60TB through events generated in other regional areas and locally at the Tier-1 at BNL. A factor of  $\sim 4$  more data was collected in the U.S. through use of parallel data stream transfer and distributed automated data management scripts than would have been possible without these grid based technologies. Ramp up is continuing for the DC2 data challenge, which will start around April 2004.

During DC1 the majority of the job execution was managed through command line scripts and Globus RSL. For the total sample more than 30,000 jobs were executed. For DC2 the job execution environment is based on Chimera/Virtual Data and the VDT. Two to three event samples of 50,000 events, each representing  $\sim 700$  jobs, have been run easily across Grid3 grid. Test samples of the order of 10s of jobs have been run interoperably across the LHC Computing Grid and Grid3. The experience from these test runs leads us to expect benefits of at least those achieved in DC1 for data collection and a decrease of a factor of  $\sim 2$  in operations manpower (increase in efficiency) of running the distributed production system.



Some References: Atlas computing workshop: August 2003

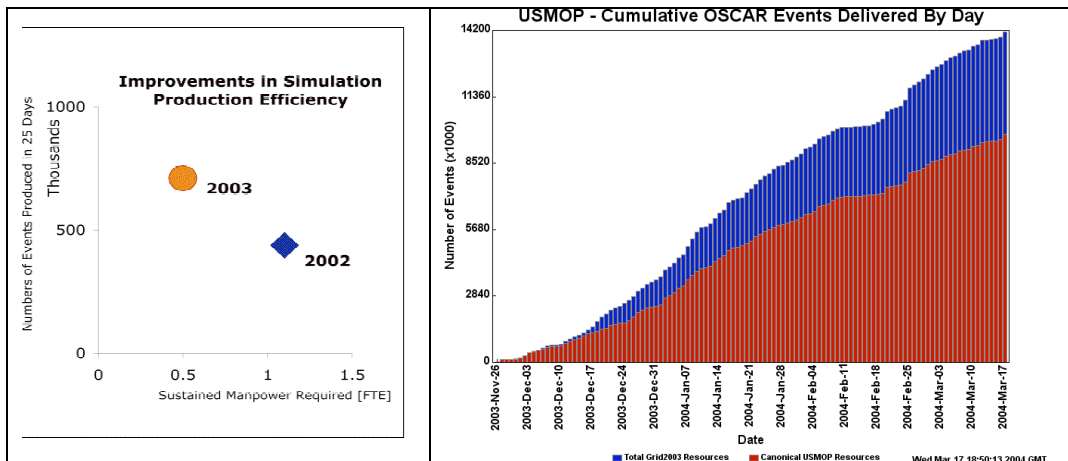
<http://wlap.physics.lsa.umich.edu/atlas/computing/workshops/2003BNL/>; MAGDA home page

<http://www.atlasgrid.bnl.gov/magda/dyShowMain.pl> ; Grid2003 US ATLAS News article

<http://www.ivdgl.org/grid2003/news/>

## US CMS

US CMS ran simulated event production on the U.S. CMS Grid testbed in 2002. At the time the 50,000 events generated reflected the largest sample of events generated on the U.S. grids. Many problems were encountered and fixed and a total sample of a million events was generated and stored using GridFTP at the Tier-1 center at Fermilab. Simulated event production in preparation for this years 2004 DC04 5% data challenge has been entirely Grid based. Jobs submitted using MOP and the CMS Distributed Production Environment based on Condor-G and Globus through use of the Virtual Data Toolkit have generated over 50 million events with an overall factor 2 more efficiency than a year or two ago. Over 75,000 jobs have been run ( cmsim: 64,000, oscar: 13,500) operated by a single FTE.

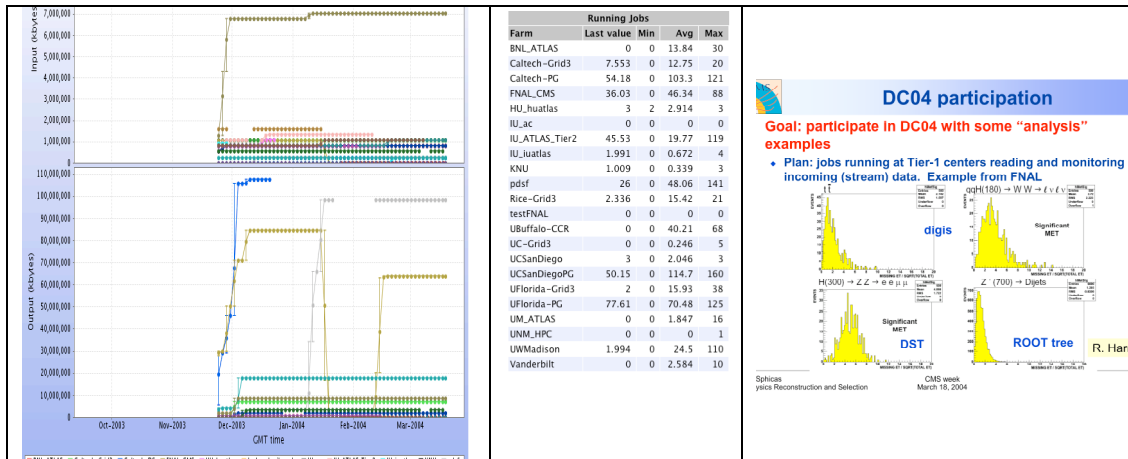


The rate of event generation was additionally increased by a factor of 50% through the opportunistic use of Grid based resources on the common grid infrastructure Grid2003. The resources of the distributed US CMS Grid are equivalent to those local to Fermilab, thus the total impact of Grid technologies is more than a factor of 2 over what would have been possible only 2 years ago

U.S. CMS currently stores more than 30TB of event in the Mass Storage System at Fermilab. More than 2/3 of that is the output from the event production, and was transferred in 1-2 GB files using GridFTP and Storage Resource Management (SRM) implementations. The efficiency and performance of the Grid based production system is measured by MonaLisa which provides system wide instrumentation and information publishing allowing the operations teams to identify changes in activity, both short term and trends, and look for problems (each point indicates the transfer of the indicated amount of data).

*KBytes transferred*

*Jobs run over 3months.....Physics from DC04*



Some References: <http://www.ppdg.net/docs/news/news-update-cmstestgrid-17may02.pdf>; MOP <http://www.uscms.org/s&c/MOP/>; CMS CMS week Mar 2004: <http://agenda.cern.ch/fullAgenda.php?ida=a041041>; monalisa <http://monalisa.cacr.caltech.edu/>;

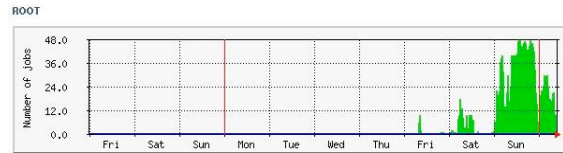
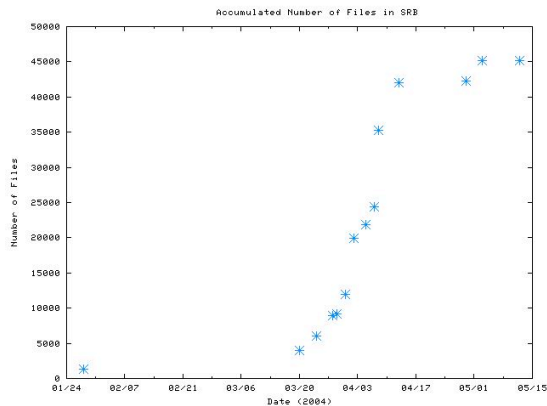
## BaBar

BaBar data management uses the extended version of the Storage Resource Broker (SRB), which supports federated MCAT catalogs for management of the meta-data in support of the automated transfer of ~10TB of data to ccin2p3 from SLAC in California.. Several thousand BaBar files (ie ROOT files) are stored in the SRB as well as 200 Conditions files. BaBar currently has ~3-5 clients using the SRB at the moment GridKA and ccin2p3 in France have used the conditions data copied via SRB for skim production. Several million events a week are being produced in this way, a number which will grow rapidly owing to the scalability of the grid-based solution.

A system for distributed analysis jobs is being piloted using simple Grid middleware [4]. Real user analysis jobs are being sent from Manchester to a mix of sites (Rutherford Laboratory, Bristol, and Manchester itself.). A small user community is developing, which is expected to grow as the spread of LCG resources makes more hardware available through this route. The BaBar experiment is making use of the Storage Resource Broker to aid in the task of distributing data to collaborating Tier A (or Tier 1) computing centers. To date practically all the experimental data in the new computing model data format (ie ROOT) have been registered in the SRB metadata catalog at SLAC. These data corresponded to over 30TB in size and number over 40K files (see figure 1). All these data are being replicate to ccin2p3, Lyon where the data are currently being analysed (see figure 2). In addition to the experiment data there are over 700 Objectivity/DB database files registered, these files contain detector condition information needed by simulation production and data processing and analysis..

The system using the SRB for distribution of conditions data for simulation production has been demonstrated at SuperComputing 2003 and this system along with the plans for data distribution using the SRB have been described in a talk at CHEP'03 [1] at SDSC, San Diego in a talk and paper at the 2003 IEEE Nuclear Science Symposium [2] and Medical Imaging Conference, Portland, Oregon and in a talk at ACAT'03 at KEK, Japan.





#### References:

- [1] <http://chep03.ucsd.edu/files/392.pdf>
- [2] <http://www.slac.stanford.edu/BFROOT/www/Computing/Offline/DataDist/DataGrids/srb/talks/ieee2003wk.pdf>
- [3] <http://www.slac.stanford.edu/BFROOT/www/Computing/Offline/DataDist/DataGrids/srb/talks/acet03ah.ppt>
- [4] BaBar Web job submission with globus authentication and afs access T Adye, R Barlow, A Forti, A McNab, S Salih, S Smith  
CHEP 2003, La Jolla

## Condor Group

The collaborative team work between the computer science groups and physics experiment teams in PPDG to deliver the benefits listed above resulted in many improvements in the grid middleware components, both in terms of capability and in robustness. These improvements are then naturally applied and made available to other applications through incorporation into the common middleware toolkits. The Virtual Data Toolkit (VDT) has been adopted by PPDG and enabled on the Alliance Testbed for evaluation.

One example is the extensions and robustness improvements in the Condor DAGMAN software, which is now benefiting the biology community at the University of Wisconsin in their execution of BLAST. They have been able to increase the number of comparisons per run from the millions to over 4 billion. For the CNS/Cyana group successful computational runs have increased from several thousand to over 25000 CPU hours.

References: DAGMAN: <http://www.cs.wisc.edu/condor/dagman/>

The Condor Exerciser application is being deployed on the Grid3 common grid environment to test the robustness and operation of the infrastructure in a systematic and ongoing basis.

## Globus Team

We choose here just one example of technology benefit and transition from the PPDG environment to another application domain. This relies, as indicated in several places, on the broad program of hardening that the PPDG approach is contributing to the deployed middleware.

GADU – the Genome Analysis and Database Update system, addresses the vital need of the emerging systems biology and DOE Genomes to Life (GTL) program for the development of high-throughput computational environments that integrate large amounts of genomic and experimental data with powerful tools and algorithms for knowledge discovery and data mining. GADU has benefited directly from several deliverables of PPDG, contributing scaling enhancements, reliability improvements, and feature development to the Grid services on which the system relies: GRAM, Condor-G, DAGMan, GridFTP, and the Replica Location Service (RLS). GADU as a whole depends on, and benefits directly from the Grid2003 infrastructure to which PPDG was a major contributor.

GADU has been used extensively by the computational biology group at Argonne National Laboratory as well as other bioinformatics organizations such as Protein Information Resources (PIR) and the Fellowship for the Interpretation of Genomes (FIG). The first GADU BLAST runs were done in March 2003, processing a peak of 59 Genomes in 24 hours. 67 CPU-days of processing time were delivered, generating 50 GB of data, using approximately 10,000 Grid jobs performing over 200,000 BLAST executions. This initial run demonstrated a greater than five-fold improvement in turnaround time: less than one hour per genome, compared to a previous average of about five hours/genome before Grid technologies were employed. GADU production runs began in August 2003, and in the first quarter of production processed 3.2 million sequences with BLAST. The first big run (for FIG) consisted of 1.8 million sequences (approximately 900MB) processed by BLAST and the result-parser, followed by subsequent monthly updates of about 80,000 sequences per month. The first production run for PIR was in November 2003, using the same process as for FIG, but on 1.2 million sequences. This was repeated in January 2004. In February 2004 we started running workflows of the BLOCKS application and the result-parser on the Grid. The initial run processed 100,000 sequences. From January-March 2004 we processed 1.3 million sequences with BLAST and 100,000 with BLOCKS.

### **Storage Resource Broker**

The Storage Resource Broker team has continued development of capabilities needed to manage distributed data collections for the high-energy physics community, and has supported application of the technology in the BaBar and CMS experiments. The SRB technology is in use in the BaBar experiment (see BaBar status report), and was used in the pre-production data challenge for CMS. The amount of data registered into SRB collections will exceed 100 TeraBytes for each experiment. The registered data is replicated between multiple countries. Because of the high wide-area-network latencies, particular attention has been paid to management of network latency through provision of bulk operations and parallel I/O stream support.

The IN2P3 team moves data from SLAC to France using two SRB servers, one located at SLAC and in Lyon. Data throughput is up to 1 TB/day (a sustained transfer rate of 11 MB/sec). This is aggregated across multiple file transfers in parallel. The data transport is highly robust, with most problems related to network infrastructure issues. A total of 30,000 files, comprising 14 TBs of data has been moved. Future plans are to import a total of 100 TBs during 2004.

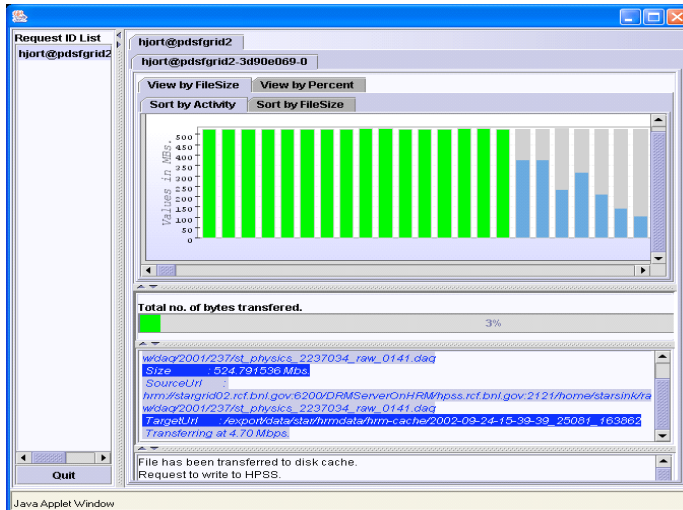
For CMS the total number of files registered in SRB is over 1 million, and the total amount of data is over 47 TBs. The data is distributed across 75 storage systems with the largest amount of data on a single resource at FNALFS (14 TBs).

### **Storage Resource Manager**

A robust replication system is now in regular use in High-Energy-Physics and Climate modeling experiments based on Storage Resource Management implementations from PPDG. Only a single command is necessary to request multi-file replication or the replication of an entire directory. A web-based tool was developed to dynamically monitor the progress of the multi-file replication process.



We have found it economical to build complex systems reusing simple robust software modules – the building block approach. SRM software modules were developed for Grid storage management. The functionality of supporting multi-file requests for clients by SRMs was applied to the file replication task. SRM uses GridFTP, to achieve secure and efficient file transfer. By using this basic file transfer service and additionally providing concurrent staging, transfer and archiving, more efficient use of the systems could be achieved.



Through the functionality of the SRM service, and developing a standard API, it has been possible for incompatible systems to interact using compatible interfaces. We have achieved robustness in the face of multiple system failures by placing monitors and recovery mechanisms in the critical paths. Since failures occur only during small percentage of the time, the recovery does not have to be efficient. In the case of file staging, transfer, and archiving failures, it is sufficient to re-issue the request.

The system developed has been in daily use by STAR, as described in that section of this document, and in frequent use by a Climate Modeling project. The modular design has permitted incremental scaling of the products to support transfer requests robustly for thousands of files, and hundreds of gigabytes of data replication in a single request. Reference: DataMover: Robust Terabyte-Scale Multi-file Replication over Wide-Area Networks SSDMB '04, [PPDG-43](#)